

Problem 1: Step Size Dilemma

a) Truncation error: $f'(x) - \frac{f(x+h) - f(x)}{h} = -\frac{f''(\xi)}{2} h$

$$\Rightarrow \left| f'(x) - \frac{f(x+h) - f(x)}{h} \right| \leq \frac{h}{2} \max_{\xi \in [x, x+h]} |f''(\xi)|$$

Round-off errors: $\varepsilon^* \stackrel{\text{def}}{=} \text{maximal relative error when subtracting}$

If f is replaced by \tilde{f} after subtraction due to finite precision:

$$\tilde{f}(x) = f(x)(1 + \varepsilon_1) \quad \text{with } |\varepsilon_1| \leq \varepsilon^*$$

$$\tilde{f}(x+h) = f(x+h)(1 + \varepsilon_2) \quad \text{with } |\varepsilon_2| \leq \varepsilon^*$$

$$\Rightarrow f'(x) - \frac{\tilde{f}(x+h) - \tilde{f}(x)}{h} = f'(x) - \frac{f(x+h) - f(x)}{h} - \frac{f(x+h)\varepsilon_2 - f(x)\varepsilon_1}{h}$$

$$= -\frac{h}{2} f''(\xi) - \frac{f(x+h)\varepsilon_2 - f(x)\varepsilon_1}{h}$$

$$\Rightarrow \left| f'(x) - \frac{\tilde{f}(x+h) - \tilde{f}(x)}{h} \right| \approx \left| -\frac{h}{2} f''(\xi) - \frac{\varepsilon_2 - \varepsilon_1}{h} f(\xi) \right|$$

use $f(x+h) \approx f(x)$

$$\leq \frac{h}{2} |f''(\xi)| + \frac{|\varepsilon_2 - \varepsilon_1|}{2} |f(\xi)|$$

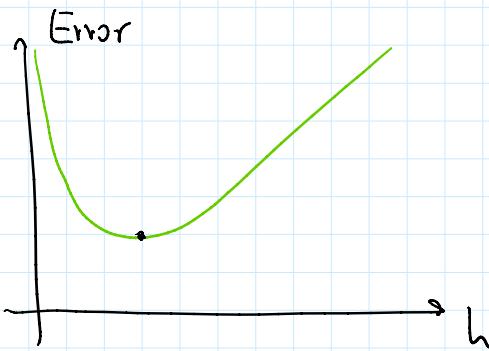
$$\leq \frac{h}{2} |f''(\xi)| + \frac{|\varepsilon_2| + |\varepsilon_1|}{h} |f(\xi)|$$

$$\leq \frac{h}{2} |f''(\xi)| + \frac{2\varepsilon^*}{h} |f(\xi)|$$

\Rightarrow total error:

$$\left| f'(x) - \frac{\tilde{f}(x+h) - \tilde{f}(x)}{h} \right| \leq \underbrace{\frac{h}{2} \max_{\xi \in [x, x+h]} |f''(\xi)|}_{\text{Truncation}} + \underbrace{\frac{2\varepsilon^*}{h} \max_{\xi \in [x, x+h]} |f(\xi)|}_{\text{roundoff}}$$

Sketch:



$$b) |S_{\text{true}}^{(d)} - \text{TD}| \leq \frac{\varepsilon |\mathcal{F}_\varepsilon| + \delta |\mathcal{F}_g|}{h^d} + |C_n| h^n$$

$$\Rightarrow 0 = -d \frac{\varepsilon |\mathcal{F}_\varepsilon| + \delta |\mathcal{F}_g|}{h^{d+1}} + n |C_n| h^{n-1}$$

$$\Leftrightarrow \frac{\varepsilon |\mathcal{F}_\varepsilon| + \delta |\mathcal{F}_g|}{h^{d+1}} = \frac{n}{d} |C_n| h^{n-1} \quad | \cdot h^{d+1}$$

$$\Leftrightarrow \varepsilon |\mathcal{F}_\varepsilon| + \delta |\mathcal{F}_g| = \frac{n}{d} |C_n| h^{n+d}$$

$$\Rightarrow h_{\text{opt}} = \left(\frac{\varepsilon |\mathcal{F}_\varepsilon| + \delta |\mathcal{F}_g|}{n/d \cdot |C_n|} \right)^{1/(n+d)}$$

$$c) \text{ Forward: } |\mathcal{F}_g| = \max_{\xi \in [x, x+h]} |f(\xi)| = |\mathcal{F}_\varepsilon|$$

$$|C_n| = \max_{\xi \in [x, x+h]} |f''(\xi)|$$

$$\Rightarrow h_{\text{opt}} = 2 \cdot \sqrt{\frac{(\varepsilon + \delta) \max_{\xi \in [x, x+h]} |f(\xi)|}{\max_{\xi \in [x, x+h]} |f''(\xi)|}}$$

$$\text{Central: } |\mathcal{F}_\varepsilon| = \max_{\xi \in [x-h, x+h]} |f(\xi)| = |\mathcal{F}_g|$$

$$|C_n| = \max_{\xi \in [x-h, x+h]} |f'''(\xi)|$$

$$\Rightarrow h_{\text{opt}} = \sqrt[3]{3(\varepsilon + \delta) \max_{\xi \in [x-h, x+h]} |f(\xi)|}$$

$$\Rightarrow h_{\text{opt}} = \sqrt[3]{\frac{3(\varepsilon + \delta) \max_{\xi \in [x-h, x+h]} |f'(\xi)|}{\max_{\xi \in [x-h, x+h]} |f'''(\xi)|}}$$

d) $f(x) = \exp(x)$ at $x=0 \Rightarrow f'(x=0) = 1$

Forward: $h_{\text{opt}} = 2 \cdot \sqrt{\frac{2 \varepsilon \exp(x)}{\exp(x)}} \approx 3 \cdot 10^{-8}$

Backward: $h_{\text{opt}} \approx 3 \cdot 10^{-8}$

Central: $h_{\text{opt}} = \sqrt[3]{\frac{6 \varepsilon \exp(x)}{\exp(x)}} \approx 9 \cdot 10^{-6}$

Comparing this to the plot last time, this was expected.

Forward and backward have the same order of accuracy whilst the central scheme is of higher order and hence achieve the optimal stepsize faster.

Problem 2: Lax-Friedrichs Method

Rewrite scheme as: $\frac{u_j^{n+1} - \frac{1}{2}(u_{j-1}^n + u_{j+1}^n) + \alpha \frac{(u_{j+1}^n - u_{j-1}^n)}{2 \Delta x}}{\Delta t} = 0$

local truncation error at (x_j, t_n) :

$$\tau_j^n = \left\{ u(x_j, t_{n+1}) - \frac{(u(x_{j-1}, t_n) + u(x_{j+1}, t_n))}{2} \right\} \cdot \frac{1}{\Delta t} + \alpha \frac{\{ u(x_{j+1}, t_n) - u(x_{j-1}, t_n) \}}{2 \Delta x} \quad (*)$$

Taylor series expansion around (x_j, t_n)

Taylor series expansion around (x_j, t_n)

$$u(x_j, t_{n+1}) = u(x_j, t_n) + \Delta t \frac{\partial u}{\partial t}(x_j, t_n) + \frac{(\Delta t)^2}{2} \frac{\partial^2 u}{\partial t^2}(x_j, t_n) + \mathcal{O}(\Delta t^3)$$

$$u(x_{j+1}, t_n) = u(x_j, t_n) + \Delta x \frac{\partial u}{\partial x}(x_j, t_n) + \frac{(\Delta x)^2}{2} \frac{\partial^2 u}{\partial x^2}(x_j, t_n) + \mathcal{O}(\Delta x^3)$$

$$u(x_{j-1}, t_n) = u(x_j, t_n) - \Delta x \frac{\partial u}{\partial x}(x_j, t_n) + \frac{(\Delta x)^2}{2} \frac{\partial^2 u}{\partial x^2}(x_j, t_n) + \mathcal{O}(\Delta x^3)$$

Insert in (*):

$$\tau_j^n = \cancel{\frac{\partial u}{\partial t}(x_j, t_n)} + \frac{\Delta t}{2} \frac{\partial^2 u}{\partial t^2}(x_j, t_n) + \mathcal{O}(\Delta t^2)$$

$$- \frac{(\Delta x)^2}{2 \Delta t} \frac{\partial^2 u}{\partial x^2}(x_j, t_n) + \frac{\mathcal{O}(\Delta x^3)}{\Delta t}$$

$$+ \alpha \cancel{\frac{\partial u}{\partial x}(x_j, t_n)} + \mathcal{O}(\Delta x^2)$$

u is exact, since $u_t + \alpha u_x = 0$. Let $\alpha = \frac{\alpha \Delta t}{\Delta x} = \text{const}$

The local truncation error then is

$$\begin{aligned} \tau_j^n &= \frac{\Delta t}{2} \frac{\partial^2 u}{\partial t^2}(x_j, t_n) - \frac{\alpha \Delta x}{2 \alpha} \frac{\partial^2 u}{\partial x^2}(x_j, t_n) + \mathcal{O}(\Delta t^2) + \mathcal{O}(\Delta x^2) \\ &= \mathcal{O}(\Delta x) + \mathcal{O}(\Delta t) \end{aligned}$$

b) Retaining error terms of the next spatial order

would solve the wave equation $u_{tt} - c^2 u_{xx} = 0$

c) See Jupyter notebook

For $\Delta x = \Delta t \ll 1$ the Lax-Tricorich's method is stable for $|\alpha| \leq 1$. In general $|\alpha| \leq 1$ for it to be stable, otherwise the solution shows strong diverging.

stable, otherwise the solution shows strong diverging oscillations.